

# Strategic Disinformation Generation and Detection

Wenxiao Yang (UC Berkeley)

Yunfei (Jesse) Yao (CUHK)

Pengxiang Zhou (USC)

# My coauthors



Wenxiao Yang (UC Berkeley)



Pengxiang Zhou (HKUST/USC)

# Research overview

- Strategic considerations of AI/ML

Disinformation detection:

Strategic Disinformation Generation and Detection (with Wenxiao Yang and Pengxiang Zhou)

Minor revision at **Management Science**

Algorithmic targeting:

Precision-Recall Tradeoff in Competitive Targeting (with Ganesh Iyer and Zemin (Zachary) Zhong)

Minor revision at **Marketing Science**

# Research overview

- Role of information in marketing strategy

- Consumer search and its marketing implications

Retargeting: A Dynamic Model of Optimal Retargeting (with J. Miguel Villas-Boas)      **Marketing Science**

Information provision: Dynamic Persuasion and Strategic Search      **Management Science**

Search fatigue: Search Fatigue, Choice Deferral, and Closure (with J. Miguel Villas-Boas)      **Marketing Science**

Advertising: Invitation to Search or Purchase? Optimal Multi-attribute Advertising      Major revision at **Management Science**

Intertemporal pricing: Non-stationary Pricing and Search (with Wee Chaimanowong)      working paper

Hidden price: A Rational Explanation of Hidden Price (with Samir Mamadehussene and Jingbo Wang)      working paper

- Privacy

Media firm's content positioning: Privacy and Polarization: An Inference-Based Framework (with Tommaso Bondi and Omid Rafieian)      **Management Science**

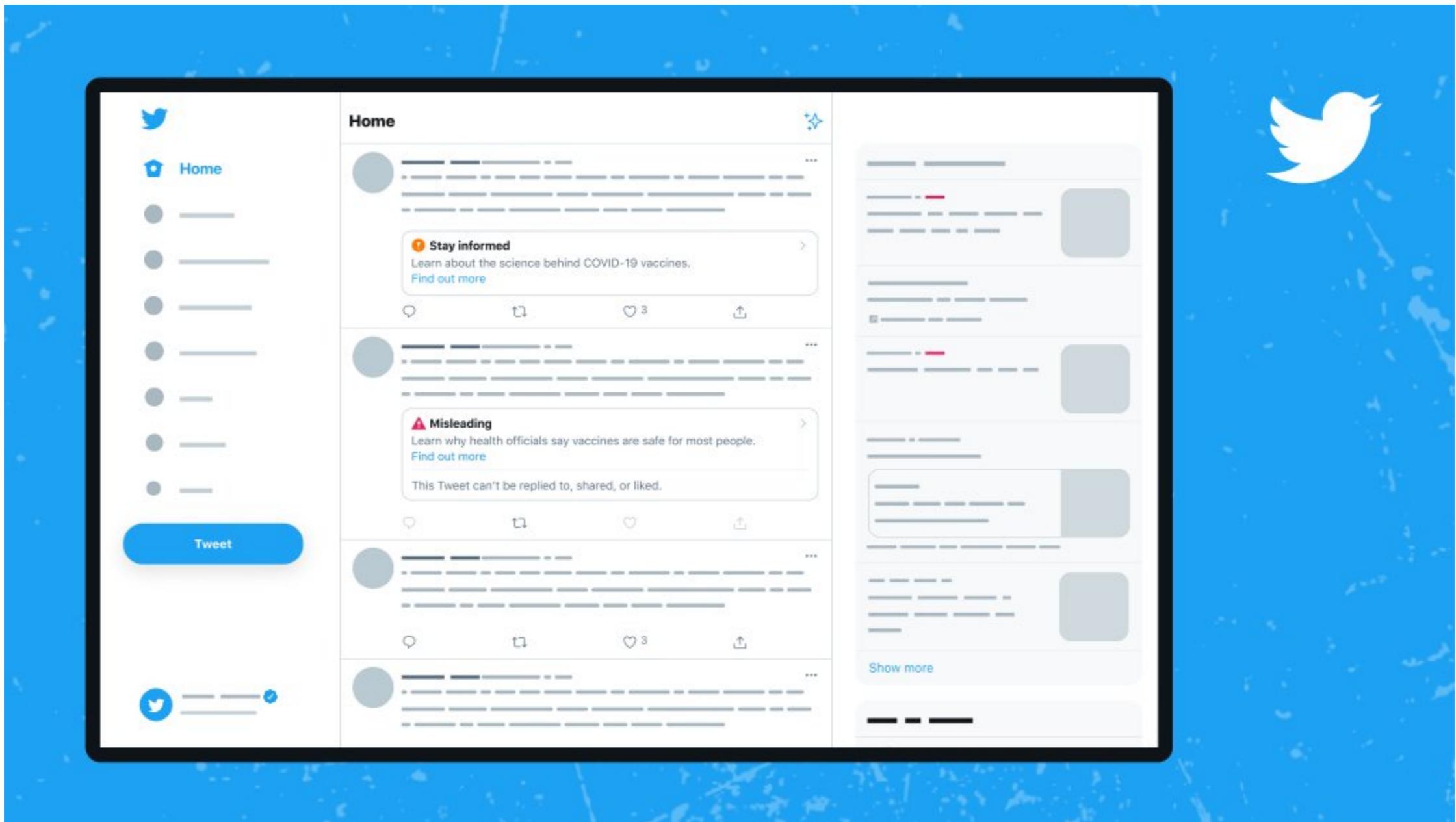
Regulatory interventions: Reputation for Privacy      **Marketing Science**

# Widespread disinformation nowadays

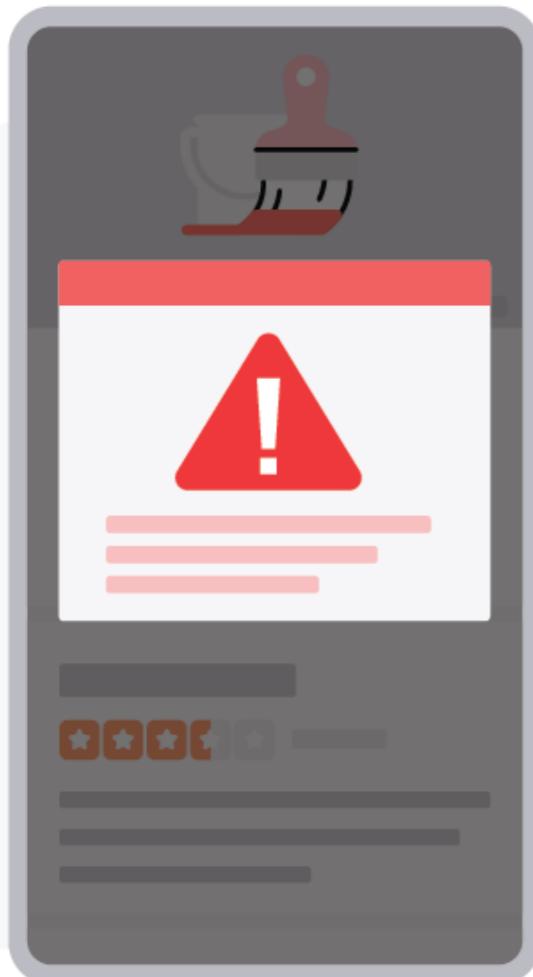
- fake reviews
- misleading posts
- fraudulent resumes
- ad fraud

Generative AI technologies further amplified such practices

# Disinformation detection (Twitter)



# Disinformation detection (Yelp)



## Suspicious Review Activity



We have noticed suspicious review activity for this business. This sort of activity can take many forms, including when a number of positive reviews originate from the same IP address or when we've identified reviews resulting from a possible **deceptive review ring**. Our **automated recommendation software** has taken this suspicious activity into account in choosing which reviews to display, but we wanted to call this to your attention because someone may be trying to artificially inflate the rating for this business.

**Got it, thanks!**

# Disinformation detection

- Platforms and regulators: deploy algorithms to detect and raise warnings about disinformation
- **Detecting** disinformation remains challenging (Callander and Wilkie, 2007; Dziuda and Salas, 2018; Mattes, Popova, and Evans, 2023)

# Dilemma in detecting disinformation

- A. increasing the likelihood of correctly recognizing deceptive content ( $\uparrow$  **true-positive rate** /  $\downarrow$  **false-negative rate**)  
cannot avoid making false-negative mistakes unless always sends an alarm
  
- B. reducing the probability of falsely identifying genuine content as deceptive ( $\downarrow$  **false-positive rate**)  
cannot avoid making false-positive mistakes unless never sends an alarm

⇒ Trade-off between

Type I error (**false-positive**) & Type II error (**false-negative**)

# Dilemma in detecting disinformation



<https://www.youtube.com/watch?v=7k1ehaE0bdU&t=1850s>

# Dilemma in detecting disinformation

- Previous work: **false negative** - a detector may fail to send an alarm when there is disinformation  
⇒ implicitly assuming that the false-positive rate is zero
- **False positives** are ubiquitous and economically significant:

J.P. Morgan views false positives as a multi-billion dollar problem;  
Global business loses more than \$100 billion annually due to false positives

# Dilemma in detecting disinformation

- Previous work: **false negative** - a detector may fail to send an alarm when there is disinformation  
⇒ implicitly assuming that the false-positive rate is zero
- **False positives** are ubiquitous and economically significant:

J.P. Morgan views false positives as a multi-billion dollar problem;

Global business loses more than \$100 billion annually due to false positives

∨

actual fraud costs

# Dilemma in detecting disinformation

- This paper: consider the possibility of **false positive** - the detector may send a false alarm without disinformation
- **Key contribution:** allow for both types of mistakes in disinformation detection.  
⇒ qualitatively different insights about strategic communication and the design of disinformation detector
- **Other main contribution:** endogenize the detector design

# Research questions

- How does the detection ability affect the incentive to generate disinformation?
- Given the practical constraints of classification technology, how should the detectors be designed?

# Related Research

- Strategic communication:
  1. verifiable disclosure (infinite cost of lying): Grossman (1981), Milgrom (1981)
  2. cheap talk (zero cost of lying): Crawford and Sobel (1982)
  3. costly lying (finite cost of lying): Kartik, Ottaviani, and Squintani (2007), Kartik (2009), Dziuda and Salas (2018), Balbuzanov (2019)

Recent work where firms and consumers interact: Villas-Boas, 2004; Guo and Zhao, 2009; Kukssov and Lin, 2010; Mayzlin and Shin, 2011; Sun, 2011; Zhang, 2013; Lauga, Ofek, and Katona, 2022; Chen, Du, and Lei, 2024; Ning, Shin, and Yu 2025.

# Related Research

- Information design: Kamenica and Gentzkow (2011), Jerath and Ren (2021); Berman, Zhao, and Zhu (2022); Ke, Lin, and Lu (2022); Pei and Mayzlin (2022); Shin and Wang (2024); Shulman and Gu (2024)
- Strategic interactions between humans and algorithms: Liang (2019), Miklos-Thal and Tucker (2019), Salant and Cherry (2020), Calvano et al. (2020), O'Connor and Wilson (2021), Montiel Olea et al. (2022), Iyer and Ke (2024), Qian and Jain (2024), Lin, Shi, and Sun (2025)

# Model

# Model

- A sender (S), a receiver (R), and a detector designer
- Receiver makes a binary decision:  $r_H$  vs .  $r_L$ 

purchasing a product from an e-commerce seller,  
re-posting social media content,  
visiting a restaurant,  
clicking on an email link,  
sending a business contact request
- Sender: high (H) type with probability  $\rho$ , low (L) type with probability  $1 - \rho$ 

Sender's private information

# Model

- Sender always wants the receiver to take action  $r_H$
- Receiver prefers to match the action with the sender's type

(sender payoff, receiver payoff)	action $r_H$	action $r_L$
type $H$ sender	$(\Delta_H^S > 0, \Delta_H^R > 0)$	$(0, 0)$
type $L$ sender	$(\Delta_L^S > 0, -\Delta_L^R < 0)$	$(0, 0)$

Table 1: Players' Payoffs

# Model

(sender payoff, receiver payoff)	action $r_H$	action $r_L$
type $H$ sender	$(\Delta_H^S > 0, \Delta_H^R > 0)$	$(0, 0)$
type $L$ sender	$(\Delta_L^S > 0, -\Delta_L^R < 0)$	$(0, 0)$

- Critical belief  $\hat{\rho}$ : receiver is indifferent between two actions  
The receiver will choose action  $r_H$  if her posterior belief exceeds the threshold  $\hat{\rho}$ ;  
choose  $r_L$  if her posterior belief is below  $\hat{\rho}$ ;  
may randomize the actions if her posterior belief is  $\hat{\rho}$ .
- Non-trivial case: receiver will take action  $r_L$  without any information,  $\rho < \hat{\rho}$

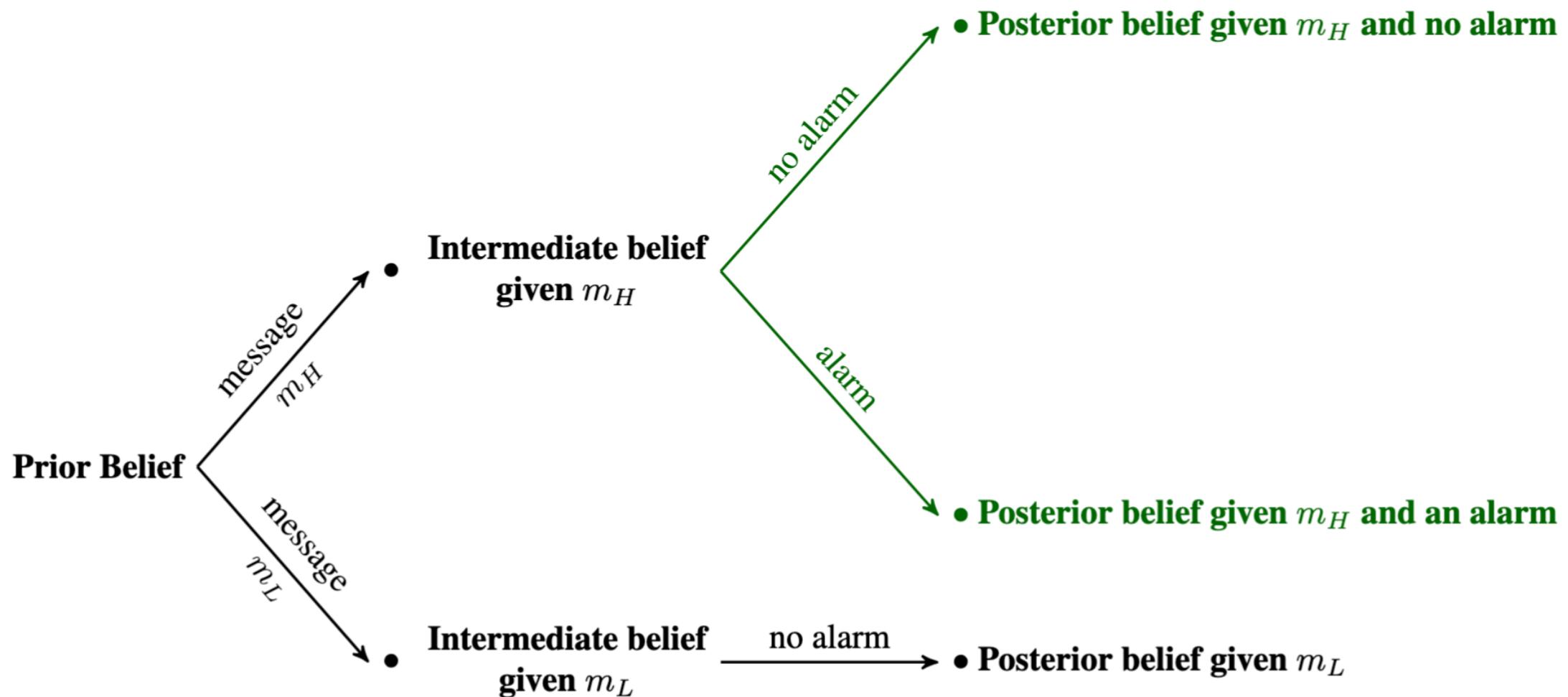
# Model

- Sender can send a message  $m \in \{m_H, m_L\}$  about his type  
non-verifiable but detectable
- Sender is **lying** if the message is not aligned with his type  
L type sends  $m_H$  or H type sends  $m_L$
- A lying cost  $C > 0$   
the sender's intrinsic aversion to lying (Gneezy 2005),  
the potential ex-post penalty for lying,  
the effort of manipulating the information
- **Lemma 1:** In equilibrium, type H sender always sends message  $m = m_H$ .

# Model

- A lie detector generates a noisy signal  $l \in \{a, na\}$  on the truthfulness of the sender's message
  - $l = a$  if the message is  $m_H$  **and** it thinks the sender is low type
  - $l = na$  otherwise
- Receiver infers the sender's type through messages from the sender and the detector.

# Receiver's belief updating



# Timing

1. The designer designs the lie detector.
2. Nature draws the sender's type  $\theta \in \{H, L\}$ .
3. The sender sends a message  $m \in \{m_H, m_L\}$  to the receiver.
4. The detector sends a signal  $l \in \{a, na\}$  to the receiver.
5. The receiver takes an action  $r \in \{r_H, r_L\}$ .

# Detector design

- A designer designs the lie detector.
- Designer's goal depends on the specific contexts
  - maximizing the receiver's payoff,
  - maximizing the high-type sender's payoff,
  - maximizing social welfare/a weighted average of sender and receiver's payoffs
  - maximizing more strategic considerations (pricing, platform entry & exit, etc.)
- **True-positive rate  $\beta$ :** probability of sending an alarm when a low-type sender mimics high type,  $Pr(l = a | m = m_H, \theta = L)$ .
- **False-positive rate  $\alpha$ :** probability of sending an alarm when the sender is high-type,  $Pr(l = a | m = m_H, \theta = H)$ .

# Detector design

- $\{\beta, \alpha\}$ : the detector's capacity (quality of detection)
- A stronger detector correctly alarms a lie more frequently and mistakenly alarms a truth-telling message less frequently

**Definition 1:** A detector  $\{\beta', \alpha'\}$  is **stronger** than a detector  $\{\beta, \alpha\}$  if and only if the following conditions hold:  
 $\beta' \geq \beta$ ,  $\alpha' \leq \alpha$ , and at least one of the inequalities is strict.

# Detector design

Receiver benefits from a low percentage of disinformation and a good detection technology.

Two channels of influencing equilibrium outcomes:

- deterring the generation of disinformation
- providing informative signals about the sender's message

# **Equilibrium concept**

Multi-stage game with incomplete information



**Perfect Bayesian Equilibrium (PBE)**

# Strategies

- Type L sender's strategy

Probability of sending message  $m_H$ :  $\sigma^S$

- Receiver's strategy

Probability of taking action  $r_H$  after seeing  $\begin{cases} m_H \text{ and no alarm: } \sigma_{na}^R \\ m_H \text{ and an alarm: } \sigma_a^R \\ m_L : \qquad \qquad \qquad \sigma_{L,na}^R \end{cases}$

# Applications

# Fake reviews and platform detection

- Review platforms such as Yelp and TripAdvisor
- Widespread phenomenon of fake reviews
  - low-quality firms invest in generating fake reviews  
⇒ mislead consumers into purchasing their products.
  - platforms implement algorithms to detect fake reviews and assist consumers by labeling potentially deceptive content
- Consumers interpret the information and make decisions given the reviews and labels

# Email marketing and alerts

- Sellers send marketing messages, hoping the recipient will click on the emails and make purchases
- Email service providers such as Outlook and Gmail develop algorithms to identify suspicious emails and issue alerts
  - ⇒ protect users from low-quality or scam content
  - design appropriate alerting rules to maximize user welfare
- Users decide whether to click on emails based on the available information

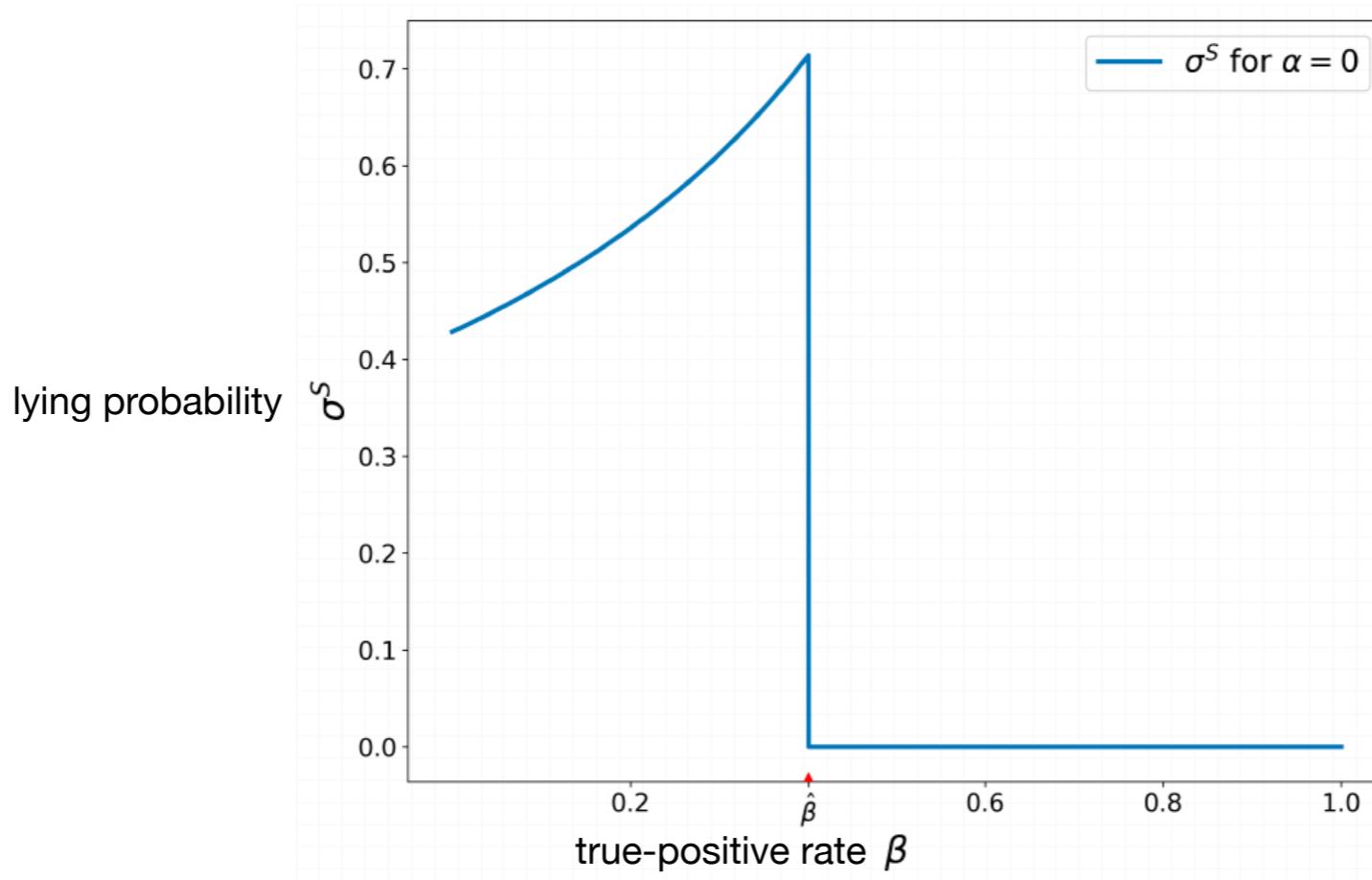
# Social media content moderation

- Social media platforms like Twitter and Facebook rely heavily on high-quality user-generated content
- Content creators may post misleading information to boost engagement metrics
  - platforms invest in content moderation to preserve user experience
  - ⇒ attach a warning label to potentially deceptive post
  - false alarms on genuine content risk alienating creators
  - failing to flag harmful content can mislead users
- Users decide whether to engage with the post based on the content and the presence or absence of label

# Benchmark

# No false-positive alarm benchmark (exogenous detector)

- If a type L sender sends message  $m_H$ , the detector sends an alarm with some probability
- False-positive rate  $\alpha = 0$ , no Type I error



$$\hat{\beta} = 1 - C/\Delta_L^S$$

# No false-positive alarm benchmark (endogenous detector)

**Lemma 4** (Endogenous detector). *The receiver's expected payoff, the high-type sender's expected payoff, and the social welfare all (weakly) increase in the true-positive rate  $\beta$ . The optimal true positive rate for the receiver is any  $\beta \geq \hat{\beta}$ . The optimal true positive rate for the high-type sender is any  $\beta \geq \min\left\{1 - \rho\Delta_H^R / [(1 - \rho)\Delta_L^R], \hat{\beta}\right\}$ . The optimal true positive rate for social welfare is any  $\beta \geq \hat{\beta}$ .*

- The more accurate the detector is, the better
- The designer always prefers a higher true-positive rate
- No trade-off !

# Main Analysis

# **Equilibrium with an exogenous detector**

# Effect of lie detection on receiver's posterior belief

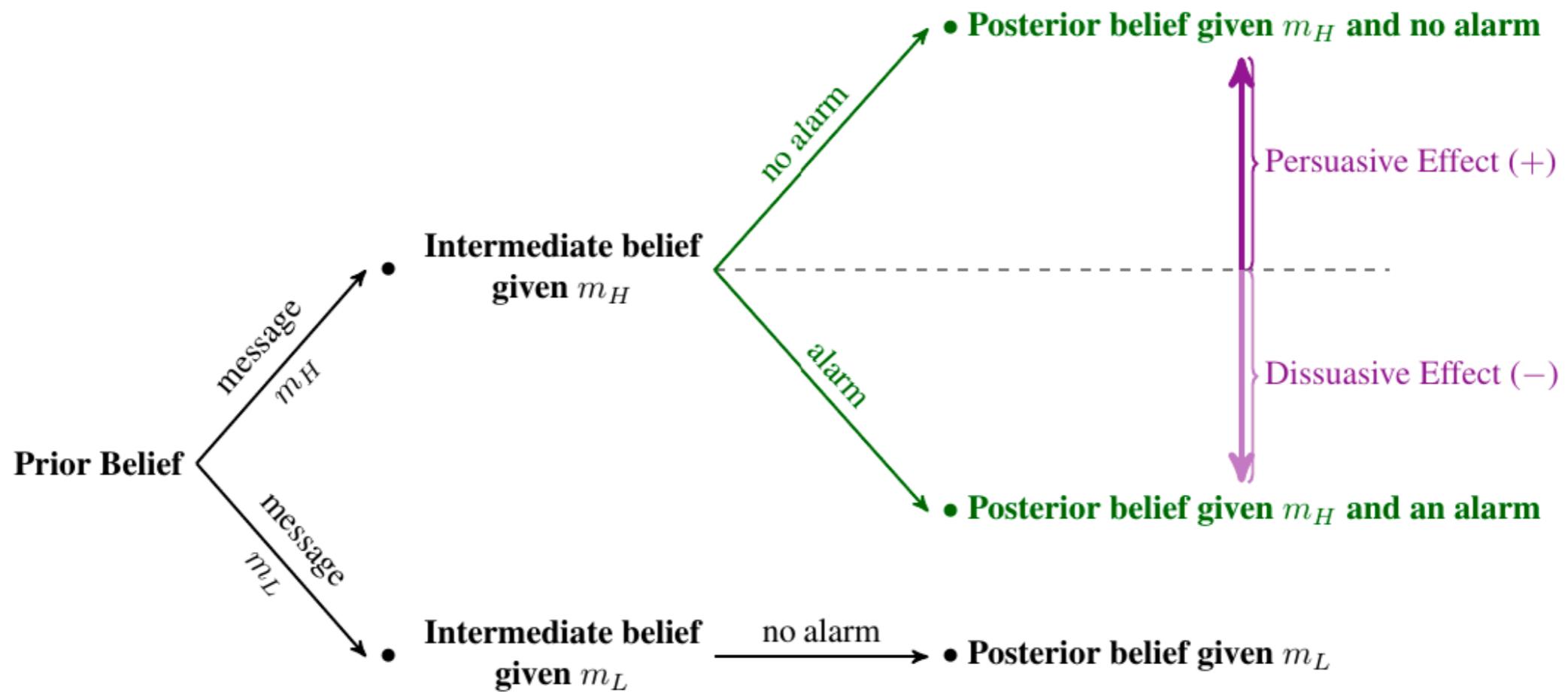


Figure 3: The Effect of Lie Detection on the Receiver's Belief.

# Persuasive effect (no alarm)

- The detector is more likely to send **no alarm** when the sender is high-type rather than low-type
- Receiver becomes more certain that the sender is high-type if she receives **no alarm**

The presence of a detector persuades the receiver to trust the sender's  $m_H$  message more

- **Persuasive effect:** the posterior belief-enhancing effect
- The persuasive effect is larger under a stronger detector

# Dissuasive effect (alarm)

- The detector is more likely to send an **alarm** when the sender is low-type rather than high-type
- Receiver becomes more certain that the sender is low-type if she receives an **alarm**

The presence of an alarm makes the receiver less trustful about the sender's  $m_H$  message

- **Dissuasive effect**: the posterior belief-reducing effect
- The dissuasive effect is larger under a stronger detector

# Dissuasive effect (alarm)

- No false-positive alarm benchmark: an alarm eliminates all the uncertainty about the sender's type

⇒ the dissuasive effect does **not** depend on  $\beta$

Two detectors with very different  $\beta$  generate the **same** effect on the posterior belief

- Main model: two detectors with the same  $\alpha$  but different  $\beta$  generate **different** dissuasive effects

⇒ variations in the dissuasive effects lead to qualitatively different equilibrium outcomes

# Effect of a stronger detector on the belief

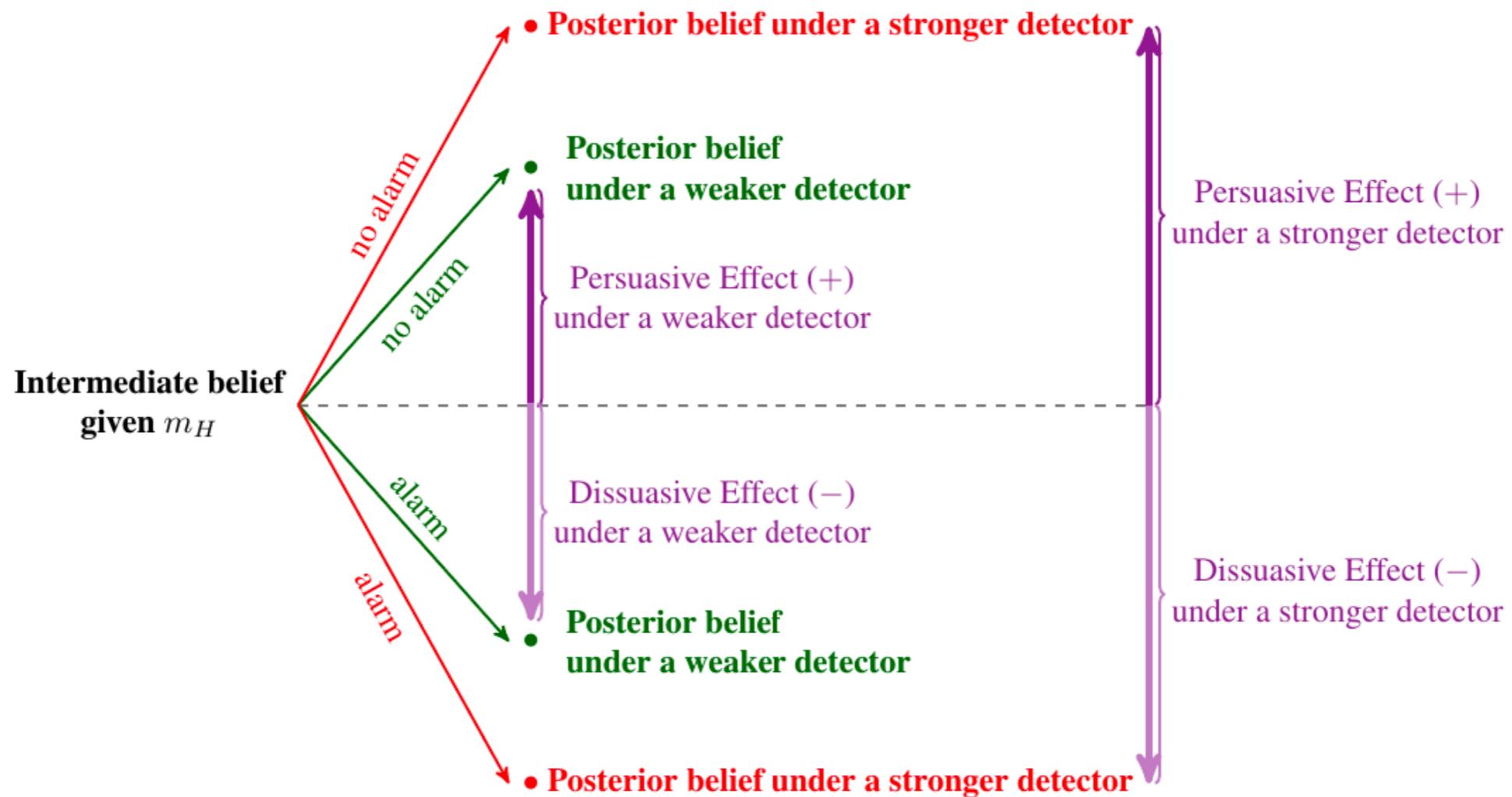
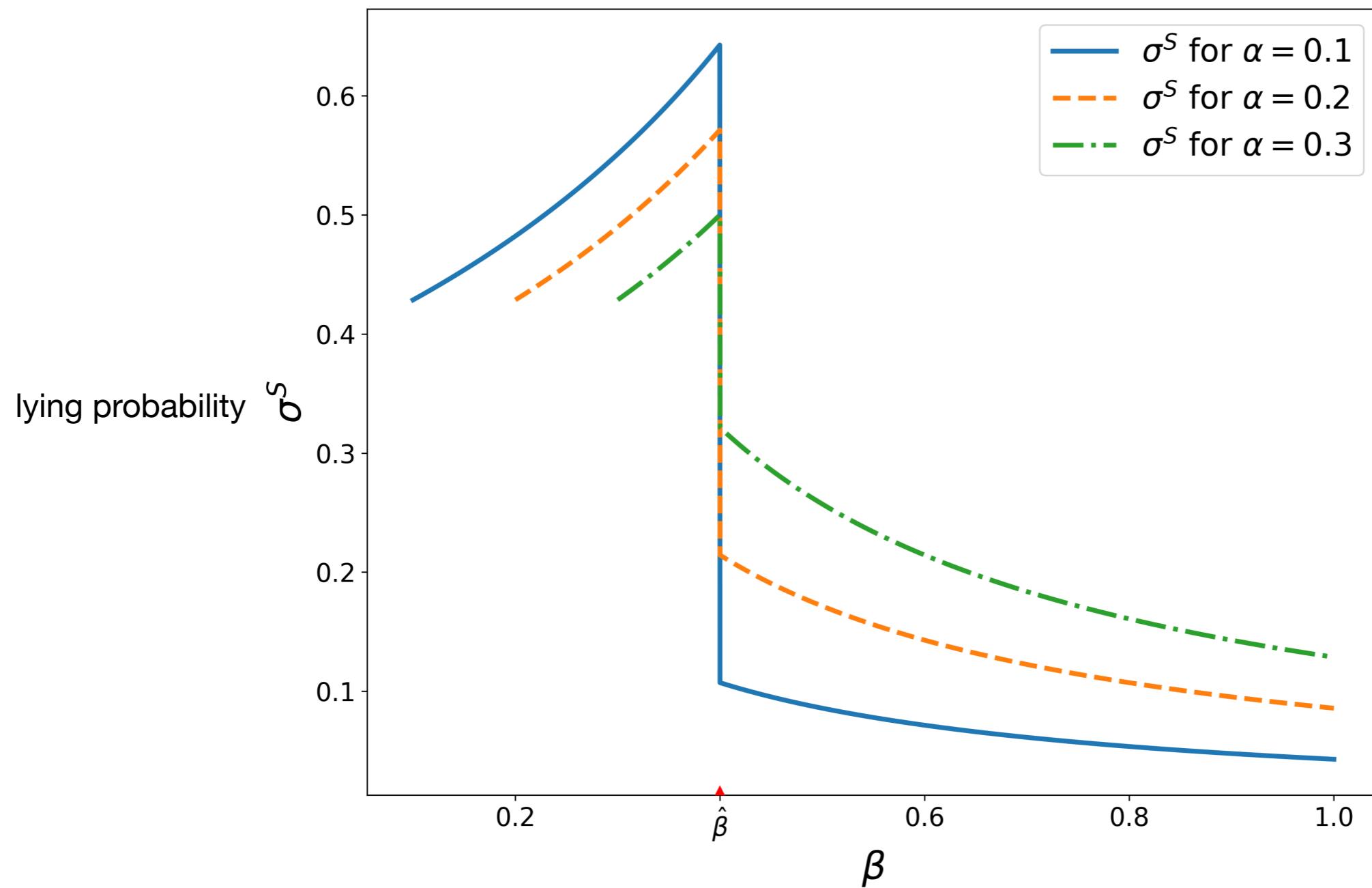
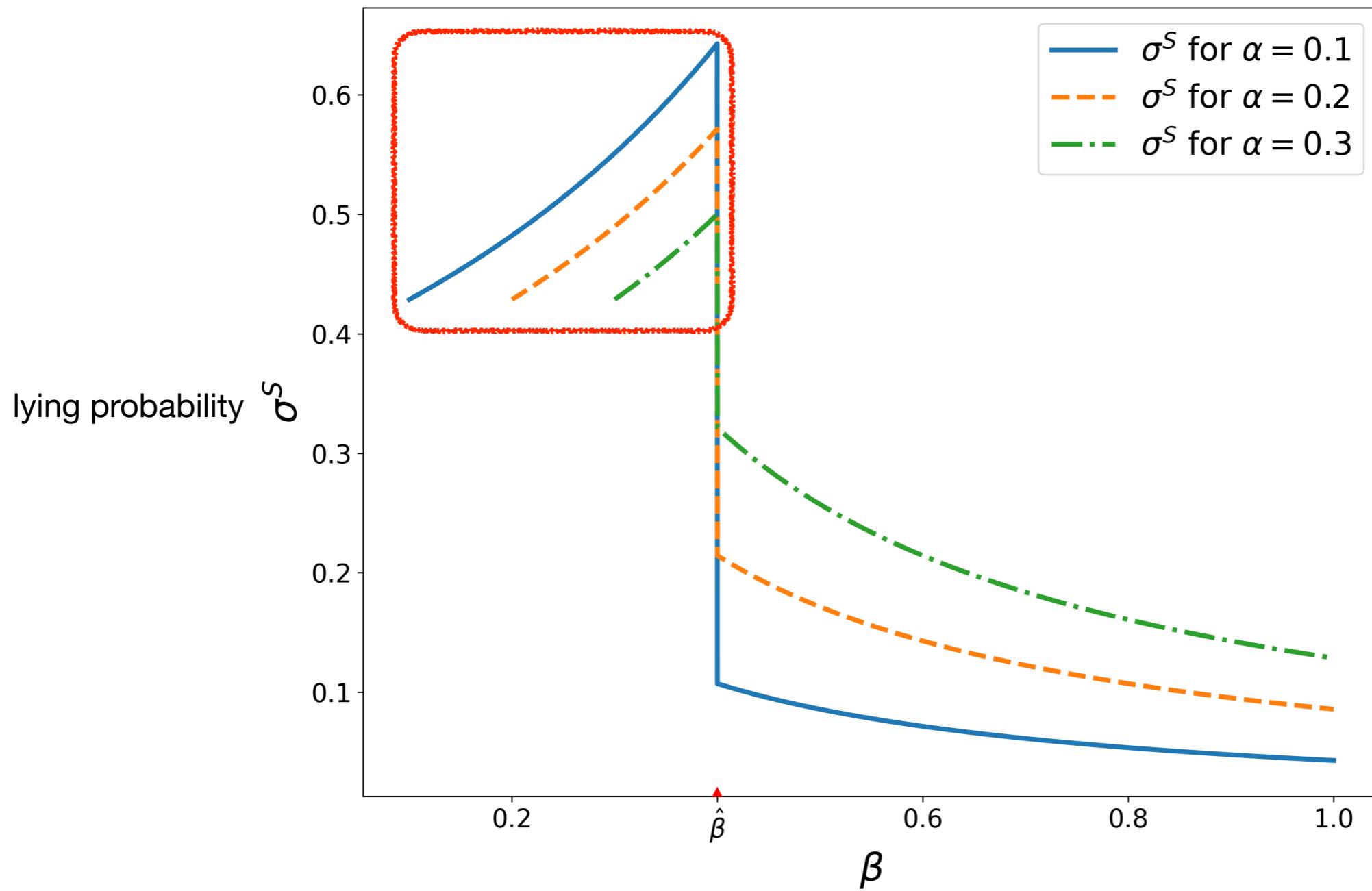


Figure 4: The Effect of a Stronger Lie Detector on the Receiver's Belief.

# Non-monotonic relationship between the detector's capacity and the sender's probability of lying



# Mechanism (low $\beta$ )



# Intuition on the non-monotonic relationship (low $\beta$ )

- low  $\beta \Rightarrow$  fail to catch many low-type senders who are lying
  - $\Rightarrow$  strong incentive for a low type sender to mimic high type
  - $\Rightarrow$  high probability of lying
  - $\Rightarrow$  low intermediate belief given message  $m_H$
  - $\Rightarrow$  low posterior belief ( $< \hat{\rho}$  given an alarm)
  - $\Rightarrow$  receiver never takes action  $r_H$  upon observing an alarm

# Intuition on the non-monotonic relationship (low $\beta$ )

- If the receiver always takes action  $r_H$  after  $m_H$  and no alarm

Low  $\beta \Rightarrow$  high benefit of lying  $(1 - \beta)\Delta_L^S >$  cost of lying  $C$

$\Rightarrow$  a low-type sender will always lie

$\Rightarrow$  low intermediate belief given message  $m_H$

$\Rightarrow$  low posterior belief also low without an alarm

$\Rightarrow$  receiver should not take action  $r_H$ , a contradiction!

# Intuition on the non-monotonic relationship (low $\beta$ )

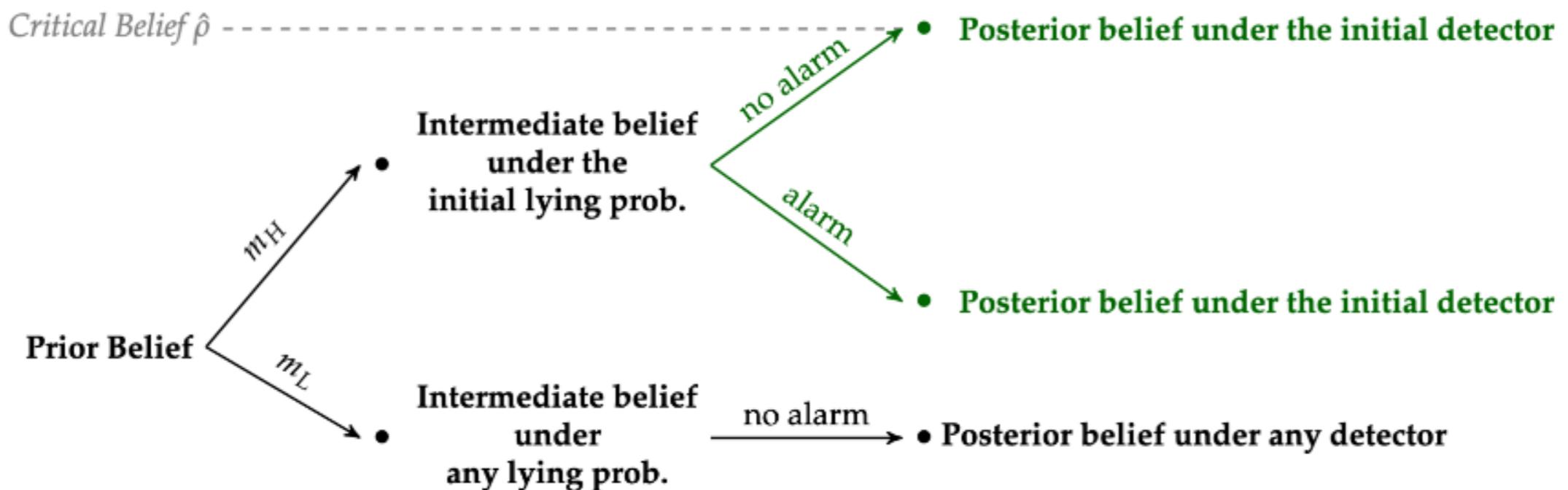
- If the receiver never takes action  $r_H$  after  $m_H$  and no alarm
  - ⇒ no low-type sender will lie
  - ⇒ only high-type senders will send message  $m_H$
  - ⇒ posterior belief = 1 given message  $m_H$  (regardless of the alarm)
    - ⇒ receiver should always take action  $r_H$  upon observing  $m_H$ , a contradiction!

# Intuition on the non-monotonic relationship (low $\beta$ )

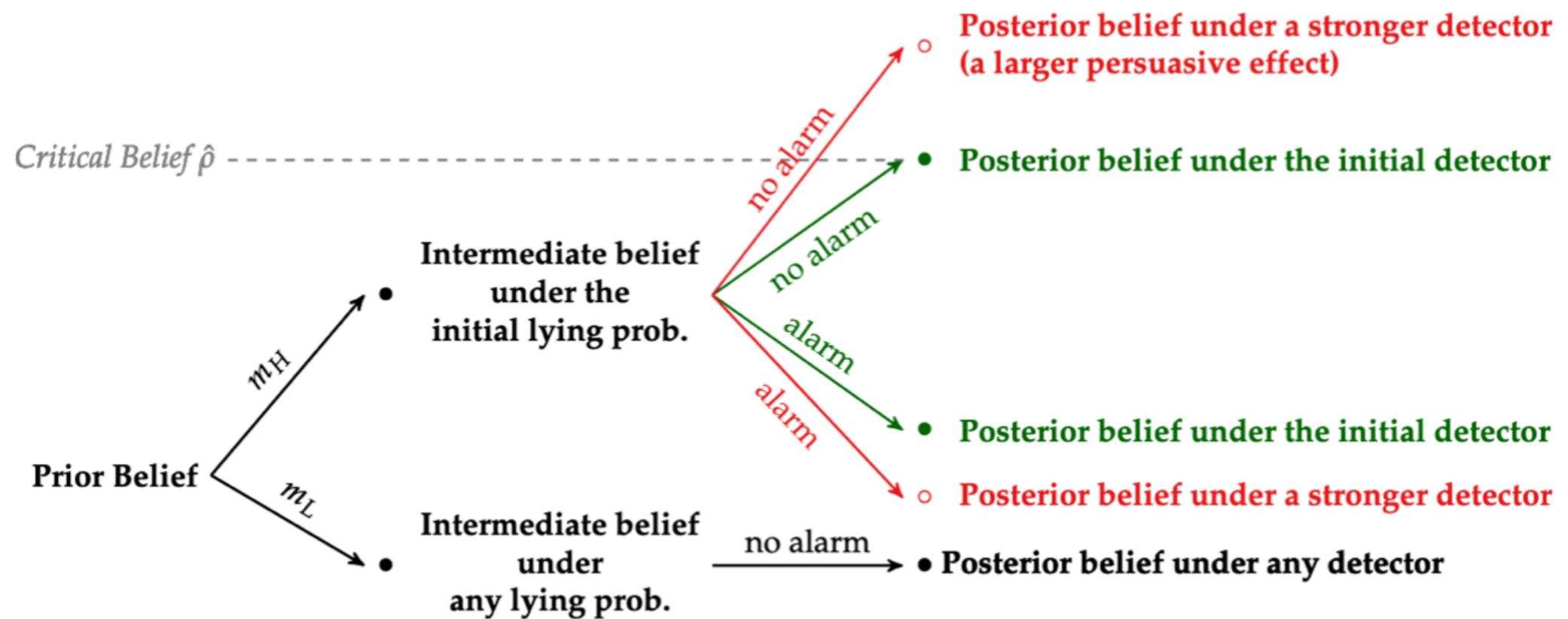
- The receiver must use a mixed strategy after  $m_H$  and no alarm

⇒ posterior belief =  $\hat{\rho}$

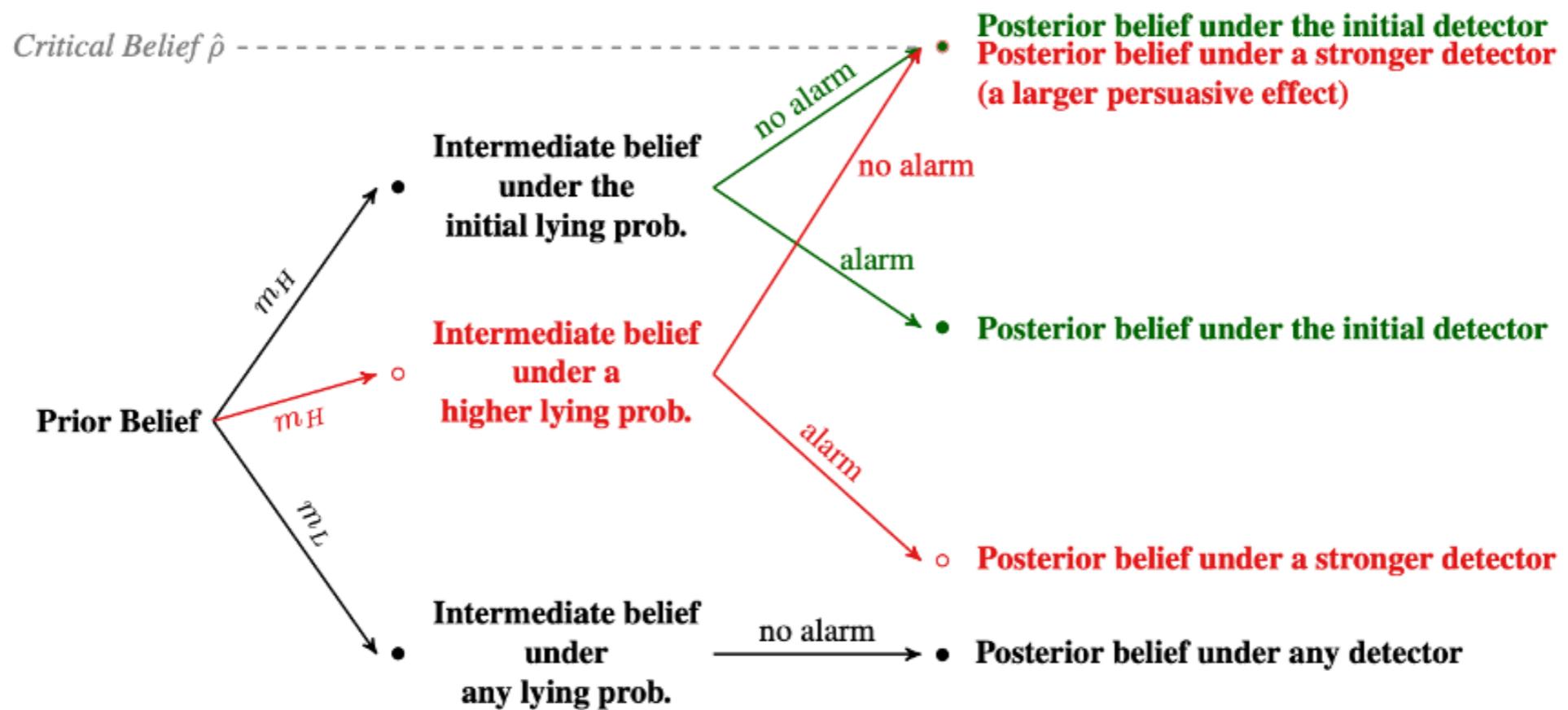
# Intuition on the non-monotonic relationship (low $\beta$ )



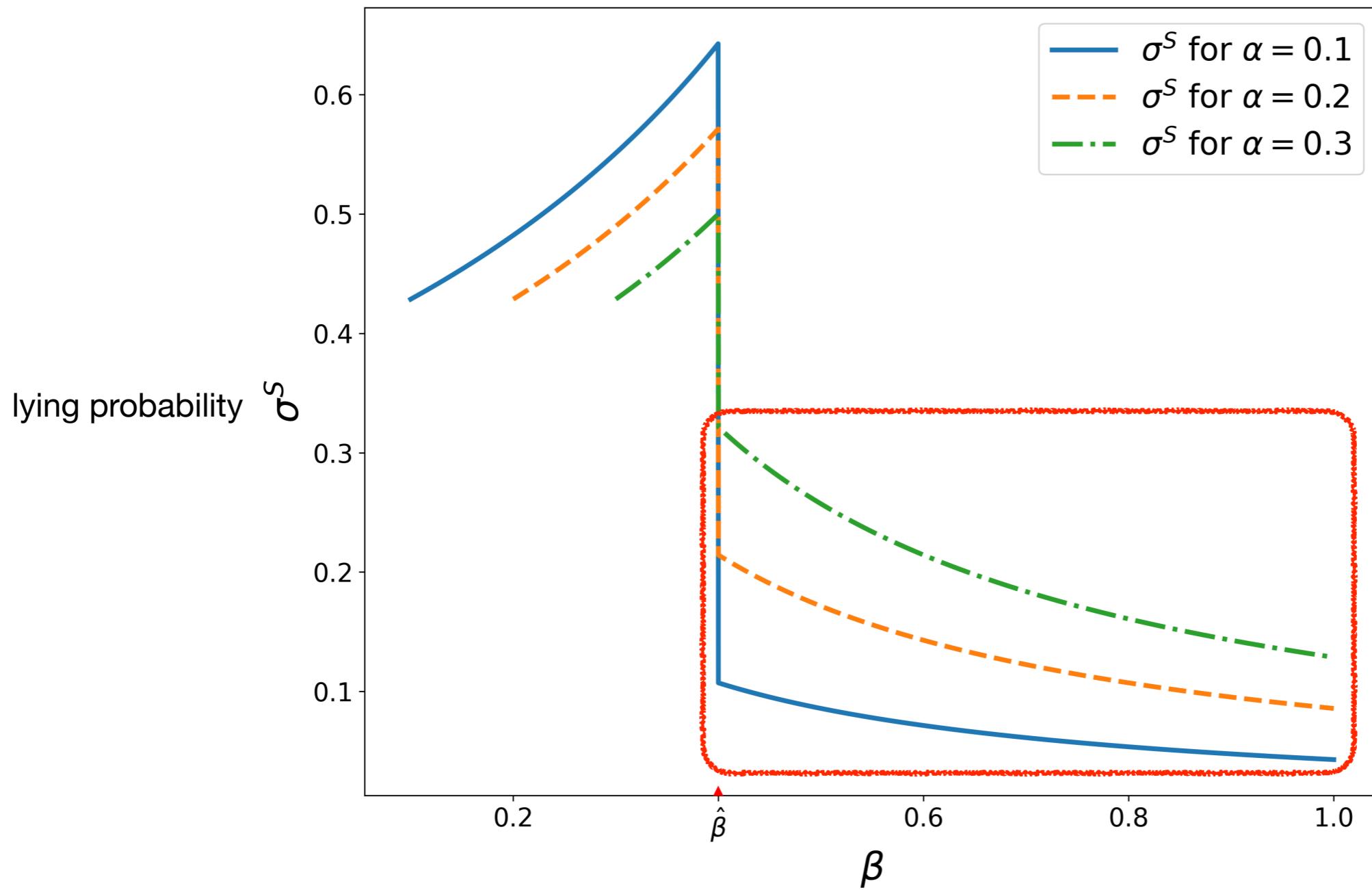
# Intuition on the non-monotonic relationship (low $\beta$ )



# Intuition on the non-monotonic relationship (low $\beta$ )



# Mechanism (high $\beta$ )



# Intuition on the non-monotonic relationship (high $\beta$ )

- High  $\beta \Rightarrow$  catch a high proportion of low-type senders who are lying
  - $\Rightarrow$  low incentive for a low type sender to mimic high type
  - $\Rightarrow$  low probability of lying
  - $\Rightarrow$  high intermediate belief given message  $m_H$
  - $\Rightarrow$  high posterior belief ( $> \hat{\rho}$  without an alarm)
  - $\Rightarrow$  receiver always takes action  $r_H$  when there is no alarm

# Intuition on the non-monotonic relationship (high $\beta$ )

- If the receiver always takes action  $r_H$  after  $m_H$  and an alarm
  - High benefit of lying  $\Delta_L^S >$  cost of lying  $C$ 
    - $\Rightarrow$  a low-type sender will always lie (pooling equilibrium)
    - $\Rightarrow$  intermediate belief = prior belief  $< \hat{\rho}$
    - $\Rightarrow$  posterior belief given an alarm  $<$  intermediate belief  $< \hat{\rho}$
    - $\Rightarrow$  receiver should not take action  $r_H$ , a contradiction!

# Intuition on the non-monotonic relationship (high $\beta$ )

- If the receiver never takes action  $r_H$  after  $m_H$  and an alarm

High  $\beta \Rightarrow$  low benefit of lying  $(1 - \beta)\Delta_L^S < \text{cost of lying } C$

$\Rightarrow$  no low-type sender will lie

$\Rightarrow$  only high-type senders will send message  $m_H$

$\Rightarrow$  posterior belief = 1 given message  $m_H$  (regardless of the alarm)

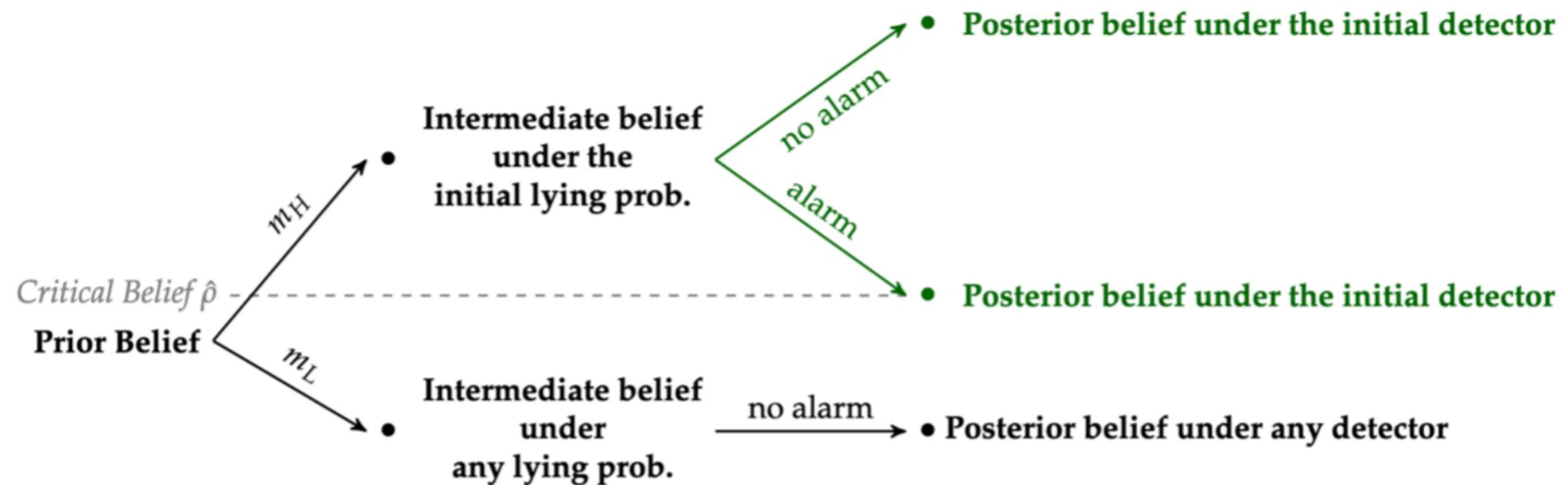
$\Rightarrow$  receiver should always take action  $r_H$  upon observing  $m_H$ , a contradiction!

# Intuition on the non-monotonic relationship (high $\beta$ )

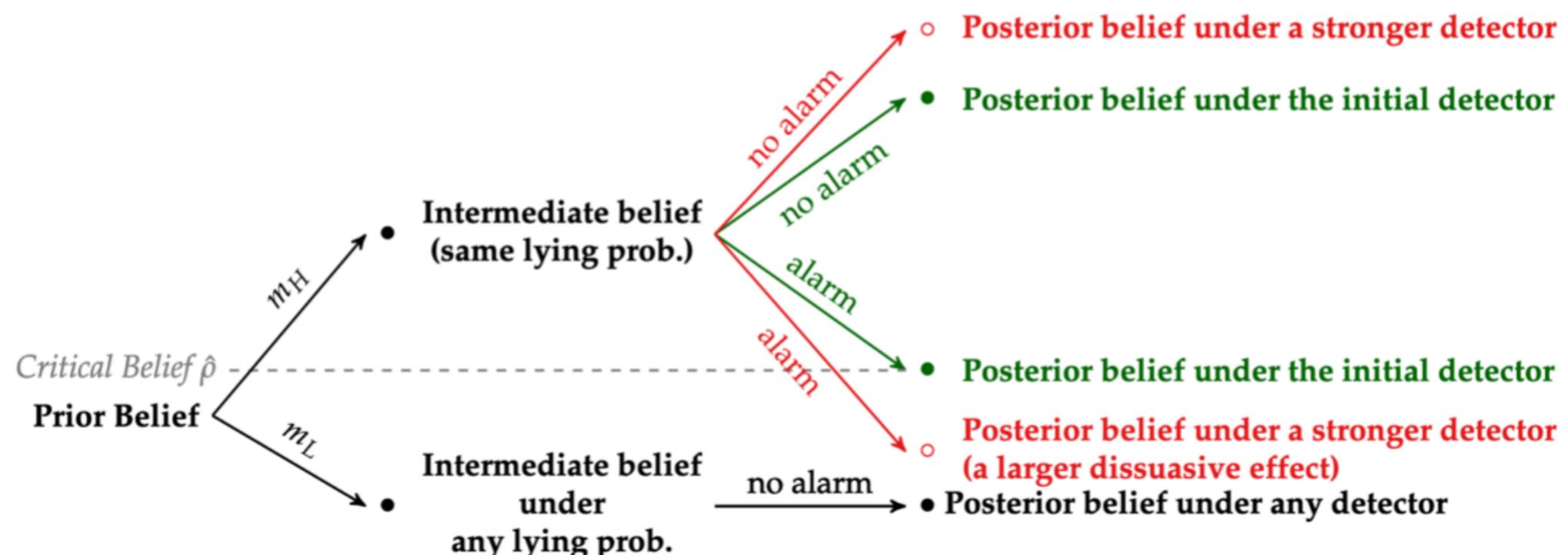
- The receiver must use a mixed strategy after  $m_H$  and an alarm

⇒ posterior belief =  $\hat{\rho}$

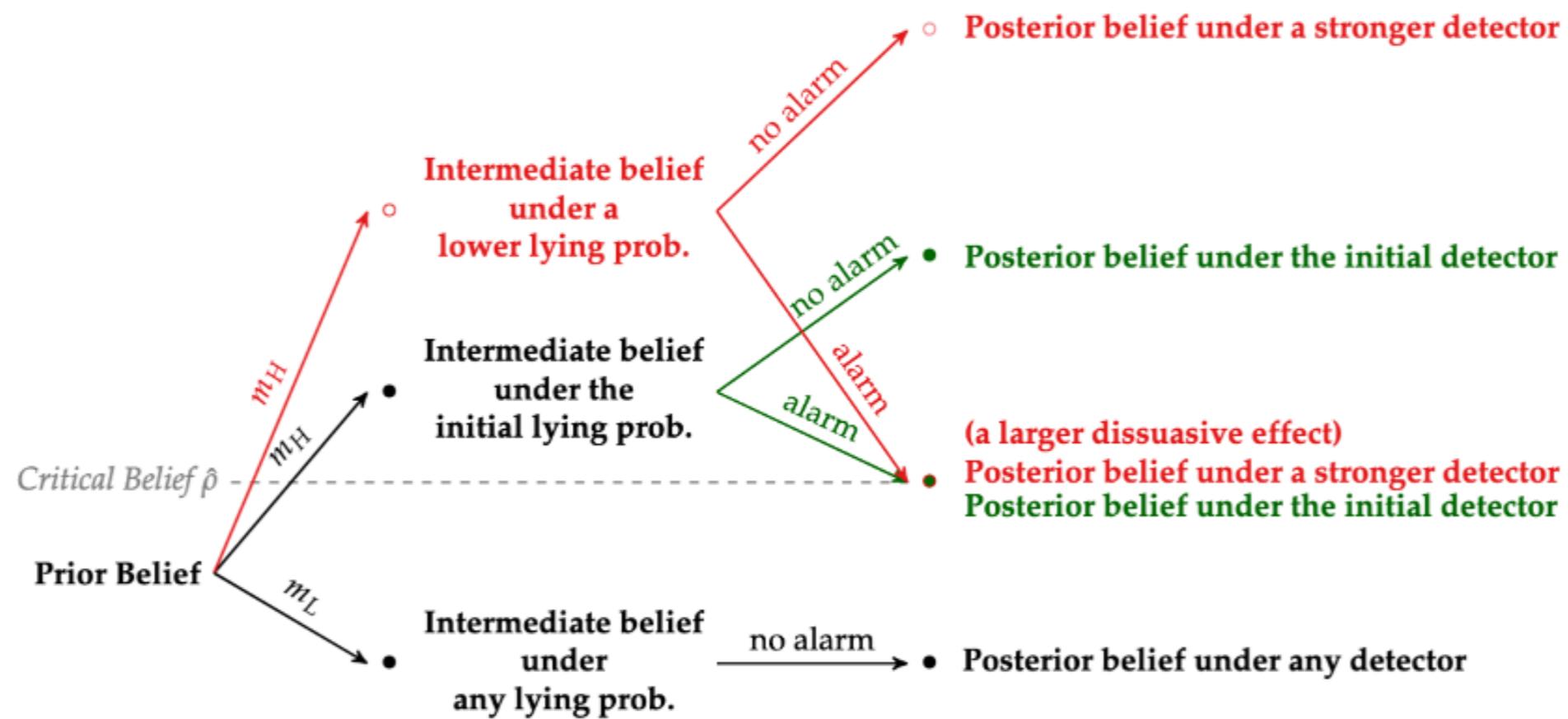
# Intuition on the non-monotonic relationship (high $\beta$ )



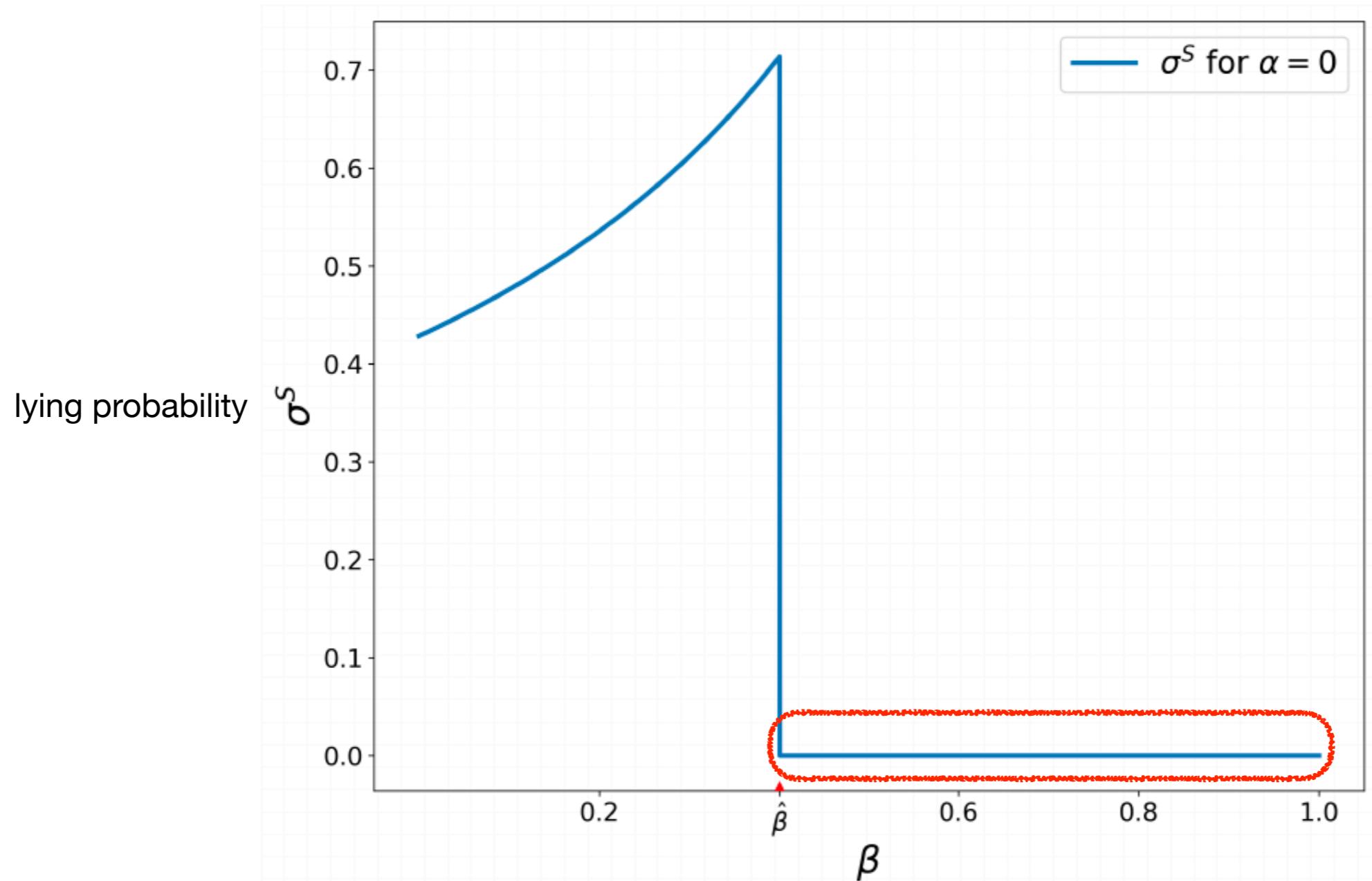
# Intuition on the non-monotonic relationship (high $\beta$ )



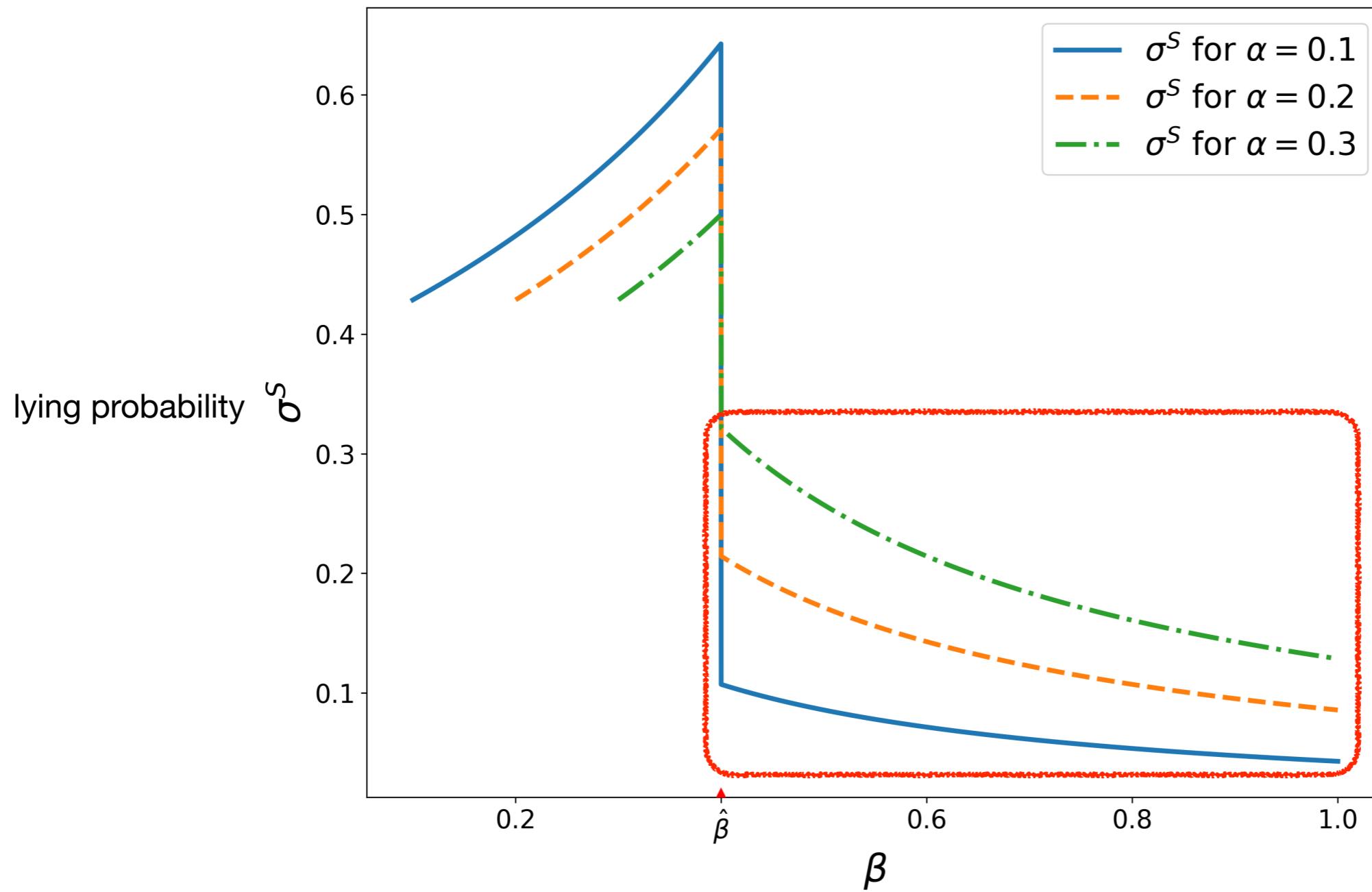
# Intuition on the non-monotonic relationship (high $\beta$ )



# No-false-positive benchmark: no disinformation under high $\beta$



# Main model: disinformation always exists



# Main model: disinformation always exists

- Suppose the sender never lies
  - ⇒ only high-type senders will send message  $m_H$
  - ⇒ an alarm must be a false-positive alarm
  - ⇒ by deviating, a low-type sender will never be caught
  - ⇒ low-type senders will lie
  - ⇒ no-lying (separating) equilibrium cannot be sustained

# **Entire Equilibrium (endogenous detector)**

# Detector design

- Designer has access to an exogenously given classifier that generates a prediction for the message's trustworthiness (technology constraint)
- Designer decides whether to send an alarm based on the prediction
- The classifier generates a binary outcome  $s \in \{s_L, s_H\}$
- $\phi(s | \theta)$ : probability of outcome  $s \in \{s_L, s_H\}$  conditional on the sender's true type  $\theta \in \{\theta_L, \theta_H\}$
- $\phi$ : the classifier's capacity (quality of classification)

Perfectly informative:  $\phi(s_H | \theta = H) = \phi(s_L | \theta = L) = 1$

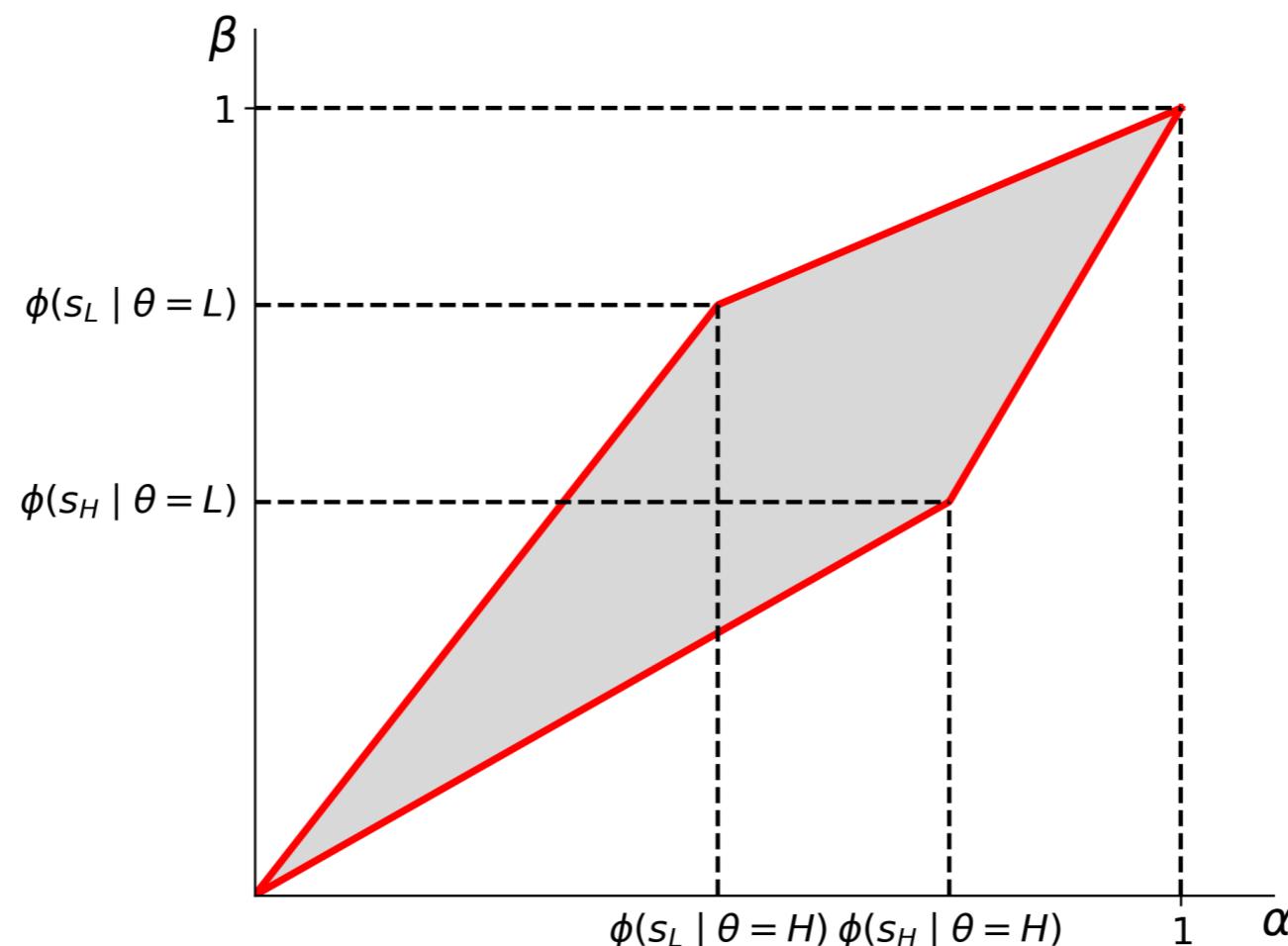
Not informative:  $\phi(s_H | \theta = H) \approx \phi(s_H | \theta = L) \& \phi(s_L | \theta = L) \approx \phi(s_L | \theta = H)$

# Detector design

- Designer's decision:
  - prob. of sending an alarm given classification outcome  $s_L$ ,  $\lambda_L = \Pr(l = a | s_L)$
  - prob. of sending an alarm given classification outcome  $s_H$ ,  $\lambda_H = \Pr(l = a | s_H)$
- $\{\lambda_L, \lambda_H\}$ : the alarm rule
- Classifier's capacity + alarm rule  $\rightarrow$  detector's capacity  $(\beta, \alpha)$

# Feasible detector space

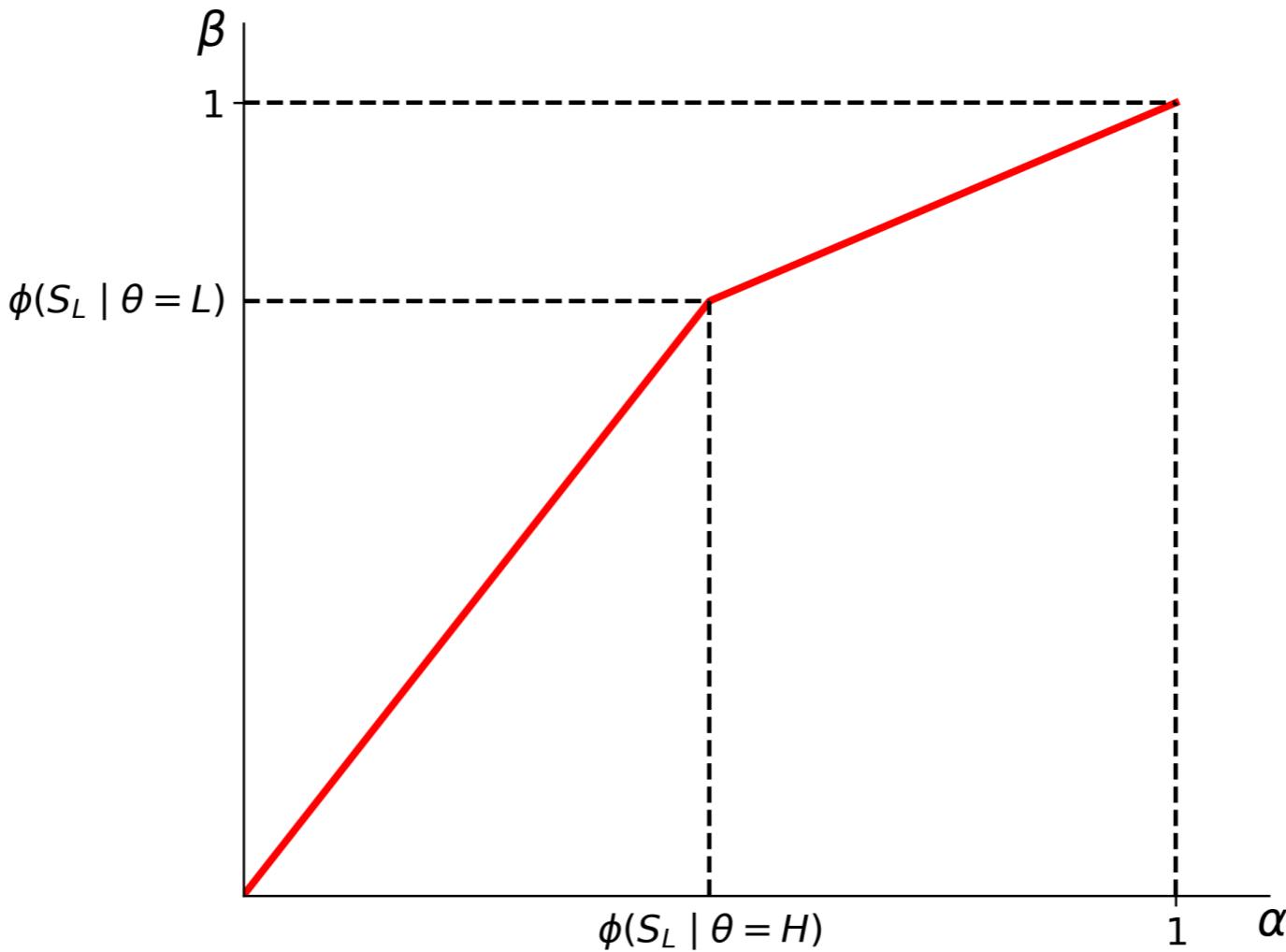
- The designer cannot obtain all detectors  $(\beta, \alpha) \in \{(\beta, \alpha) \mid 0 \leq \alpha < \beta \leq 1\}$  due to the constraint of the classifier's capacity



# Effect of lie detection on payoffs

- For a given  $\alpha$ , a higher  $\beta$  benefits the receiver and high-type sender, whereas hurts the low-type sender
  - For a given  $\beta$ , a lower  $\alpha$  makes all players better off
- ⇒ The designer always chooses the **lowest feasible false-positive rate**  $\alpha$  given any true-positive rate  $\beta$ .

# Receiver operating characteristic (ROC) curve



- Pareto frontier of the classification outcome
- the lowest feasible false-positive rate  $\alpha$  given any true-positive rate  $\beta$

# Optimal false-positive rate and alarm Rule given true-positive rate

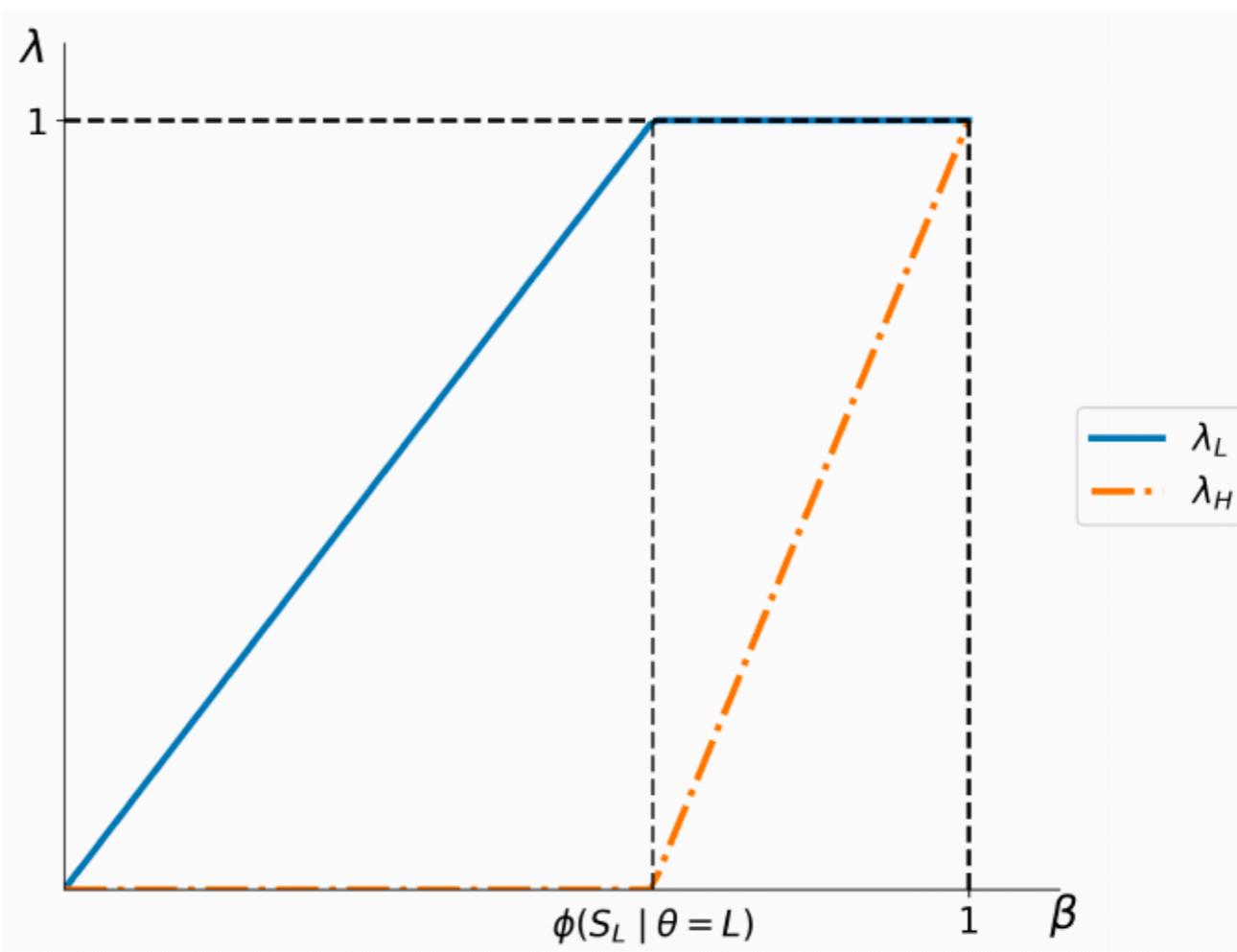
**Lemma 5** (Optimal False-positive Rate and Alarm Rule Given True-positive Rate). *For a given true-positive rate  $\beta$ , the detector's optimal false-positive rate, denoted by  $\alpha^*(\beta; \phi)$ , is*

$$\alpha^*(\beta; \phi) = \begin{cases} \frac{\phi(s_L|\theta=H)}{\phi(s_L|\theta=L)}\beta, & \text{if } \beta \leq \phi(s_L|\theta=L) \\ \frac{\phi(s_H|\theta=H)}{\phi(s_H|\theta=L)}\beta + 1 - \frac{\phi(s_H|\theta=H)}{\phi(s_H|\theta=L)}, & \text{if } \beta > \phi(s_L|\theta=L), \end{cases}$$

which increases in  $\beta$ . The detector  $\{\beta, \alpha^*(\beta; \phi)\}$  can be achieved by the alarm rule

$$\lambda_L^*(\beta) = \begin{cases} \frac{\beta}{\phi(s_L|\theta=L)}, & \text{if } \beta \leq \phi(s_L|\theta=L) \\ 1, & \text{if } \beta > \phi(s_L|\theta=L) \end{cases}, \quad \lambda_H^*(\beta) = \begin{cases} 0, & \text{if } \beta \leq \phi(s_L|\theta=L) \\ \frac{\beta - \phi(s_L|\theta=L)}{\phi(s_H|\theta=L)}, & \text{if } \beta > \phi(s_L|\theta=L). \end{cases}$$

# Optimal alarm rule given a true-positive rate



- The designer wants to achieve a given true-positive rate while minimizing the false-positive rate.

# Optimal Design of Detector

- Optimal choice of a feasible detector,  $\{\beta^*, \alpha^*\}$
- Previous lemma: optimal false-positive rate is  $\alpha^*(\beta; \phi)$  given any true-positive rate  $\beta$ .  
⇒ only need to pin down the optimal true-positive rate  $\beta^*$
- **Lemma 6:** Sequential derivation  $\Leftrightarrow$  simultaneous optimization

# High/low capacity of a classifier

- The classifier can better distinguish the sender's type if it generates  $s_H$  more frequently when the sender is high type than low type;
- It also can better distinguish the sender's type if it generates  $s_L$  more frequently when the sender is low type than high type.

**Definition 2.** A classifier has a high capacity if  $\phi(s_H|\theta = H)/\phi(s_H|\theta = L) \geq (1 - \rho)\Delta_L^R/(\rho\Delta_H^R)$  and  $\phi(s_L|\theta = L)/\phi(s_L|\theta = H) \geq (\Delta_L^S - C)\rho\Delta_H^R/[\Delta_L^S\rho\Delta_H^R - (1 - \rho)\Delta_L^RC]$ . Otherwise, it has a low capacity.

# Maximizing receiver's payoff

- **Proposition 3:** The  $\beta^*$  that maximizes the receiver's expected payoff is:

(1) If the classifier has a low capacity:

any  $\beta \in [\hat{\beta}, \max\{\hat{\beta}, \phi(s_L | \theta = L)\}]$ , which minimizes the low-type sender's equilibrium probability of lying,

(2) If the classifier has a high capacity:

any  $\beta \in [\hat{\beta}, \max\{\hat{\beta}, \phi(s_L | \theta = L)\}]$  if the lying cost is high;

$\beta = \phi(s_L | \theta = L)$ , if the lying cost is low.

# Maximizing receiver's payoff

Two channels of influencing equilibrium outcomes:

- deterring the generation of disinformation
  - High lying cost: easy to discourage low-type senders from lying
  - Low-capacity classifier: hard to provide highly informative alarm signals  
⇒  $\beta^*$  minimizes the low-type sender's probability of lying (discrete jump at  $\hat{\beta}$ )
- providing informative signals about the sender's message
  - Low lying cost/high-capacity classifier: takes full advantage of the region where a unit increase in  $\beta$  leads to a small increase in  $\alpha$

# Maximizing high-type sender's payoff

- **Proposition 4:** The  $\beta^*$  that maximizes the high-type sender's expected payoff is:
  - (1) If the classifier has a low capacity: any  $\beta \in [\hat{\beta}, \max\{\hat{\beta}, \phi(s_L | \theta = L)\}]$ ;
  - (2) If the classifier has a high capacity:  
$$\beta_1 := [\rho \Delta_H^R - (1 - \rho) \Delta_L^R] / [\rho \Delta_H^R \phi(s_L | \theta = H) / \phi(s_L | \theta = L) - (1 - \rho) \Delta_L^R] < \hat{\beta}$$
, which is decreasing in  $\phi(s_L | \theta = L) / \phi(s_L | \theta = H)$ .

# Maximizing high-type sender's payoff

- The sender prefers the receiver to always take action  $r_H$  conditional on message  $m_H$  and no alarm (large enough  $\beta$ )
- Fixing a  $\beta$ , the detector has a lower  $\alpha$  if the classifier has a higher capacity.
  - ⇒ Larger persuasive effect
  - ⇒ Higher posterior belief
  - ⇒ detector can induce action  $r_H$  even if  $\beta$  is adjusted downward
- The designer has no incentive to further increase  $\beta$  (leads to more  $\alpha$ )
- When the classifier has a high capacity, *counter-intuitively*, the optimal detector alarms a **smaller** percentage of disinformation when its underlying classifier is **better** at distinguishing the sender's type.

# Maximizing

$$\mathbb{E} U^R(\beta) + w_H \rho \mathbb{E} U_H^S(\beta) + w_L (1 - \rho) \mathbb{E} U_L^S(\beta)$$

**Proposition 5.** *If the classifier has a low capacity, any  $\beta \in [\hat{\beta}, \max\{\hat{\beta}, \phi(s_L | \theta = L)\}]$  is optimal. If the classifier has a high capacity, the set of optimal true-positive rates is:*

$$\begin{cases} \mathcal{B}(w_H, w_L) & \text{if } C < \tilde{C}(w_H, w_L) \\ \mathcal{B}(w_H, w_L) \cup [\hat{\beta}, \max\{\hat{\beta}, \phi(s_L | \theta = L)\}] & \text{if } C = \tilde{C}(w_H, w_L) \\ [\hat{\beta}, \max\{\hat{\beta}, \phi(s_L | \theta = L)\}] & \text{if } C \in (\tilde{C}(w_H, w_L), \Delta_L^S(1 - \beta_1)] \end{cases},$$

where

$$\mathcal{B}(w_H, w_L) := \begin{cases} \{\phi(s_L | \theta = L)\} & \text{if } n_0 w_H + l_0 w_L < 1 \\ [\beta_1, \phi(s_L | \theta = L)] & \text{if } n_0 w_H + l_0 w_L = 1 \\ \{\beta_1\} & \text{if } n_0 w_H + l_0 w_L > 1 \end{cases},$$

$\tilde{C}(w_H, w_L)$  is a continuous function increasing in both  $w_H$  and  $w_L$ , and  $n_0$  and  $l_0$  are strictly positive constants.

# Maximizing

$$\mathbb{E} U^R(\beta) + w_H \rho \mathbb{E} U_H^S(\beta) + w_L (1 - \rho) \mathbb{E} U_L^S(\beta)$$

- If  $w_H = w_L = 1$ : maximizing social welfare
  - Low-capacity classifier: senders' and receivers' incentives are aligned
  - High-capacity classifier: the sender prefers a lower  $\beta$  than the receiver
- ⇒ optimal detector is a compromise between their preferences
- ⇒ greater weight on senders,  $\uparrow w_H, w_L \rightarrow$  lower  $\beta^*$

# Comparison with the no false-positive alarm benchmark

- No false-positive alarms: no trade-off in the detector design  $\Rightarrow$  designer always prefers a higher  $\beta$ .
- With false-positive alarms: designer strictly prefers intermediate  $\beta$ .

False-positive alarms + players' strategic responses  $\rightarrow$  a higher  $\beta$  may reduce the receiver's payoff, the high-type sender's payoff, and social welfare

# **Extensions**

# Extensions

- Restriction on the alarm rule:  $\lambda_H = 0$
- Endogenous commission fee and strategic platform entry

# Restriction on the alarm rule:

$$\lambda_H = 0$$

- Detector never sends an alarm when the classifier predicts signal
  - High lying cost: results and underlying mechanisms are similar to the main model
  - Low lying cost: low-type senders have a strong incentive to lie
- ⇒ designer prefers a detector with a high true-positive rate to reduce the lying probability
- ⇒ is constrained by how much to increase the true-positive rate

# Endogenous commission fee and strategic platform entry

- Platform's tool: information → Pricing + information
  1. Platform first decides the commission rate and detector design
  2. High- and low- quality sellers simultaneously decide whether to enter the platform
  3. Entry decisions are revealed, and each seller who entered chooses the message
  4. The detector sends a signal for each message
  5. Consumers decide whether to purchase a product from each seller

# Endogenous commission fee and strategic platform entry

- If both types of seller enter, the properties of the optimal detector are consistent with those in the main model

Platform pricing (commission fee) is a ***strategic complement*** to detector design

Condition: low lying cost and high classifiers capacity

Intuition: by lowering commission rate, the platform

+: gain additional revenue from low-quality sellers

-: sacrifice profits from high-quality firms

Under the condition, it can induce broader seller participation with a smaller discount on commission rate (higher benefit & lower cost)

- If the commission rate is sufficiently high → only high-quality sellers enter

Platform pricing (commission fee) is a ***strategic substitute*** for detector design

# Descriptive value

- Qualitatively different insights about the relationship between the sender's probability of lying and the detector's accuracy when allowing false-positive alarms
- Without false-positive alarms: two detectors with different true-positive rates generate the same dissuasive effect
- With false-positive alarms: two detectors with the **same** false-positive rate but different true-positive rates generate different dissuasive effects
- Variations in the dissuasive effects  $\Rightarrow$  non-monotonic relationship between the sender's probability of lying and detector's accuracy

# Prescriptive value

- Characterize the optimal design of the detector in the presence of practical limitations
- Possibility of false-positive alarms  $\Rightarrow$  the designer should not choose the largest true-positive rate
- The optimal detector may raise alarms about a smaller percentage of disinformation when its underlying classifier is better at distinguishing the sender's type
- Qualitatively different and counter-intuitive findings  $\Rightarrow$  importance of considering the interaction between senders' strategic behavior and both types of mistakes by the detection technology in practice

# **Thanks!**